

**This Page is Inserted by IFW Indexing and Scanning  
Operations and is not part of the Official Record**

**BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☒ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** \_\_\_\_\_

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.**

(19)



JAPANESE PATENT OFFICE

## PATENT ABSTRACTS OF JAPAN

(11) Publication number: **08211895 A**(43) Date of publication of application: **20 . 08 . 96**

(51) Int. Cl.

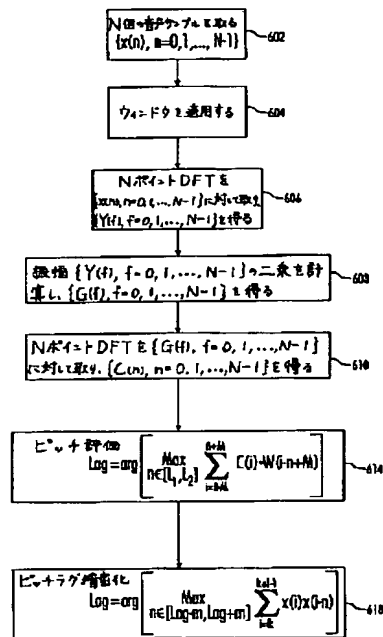
**G10L 3/02****G10L 7/04****G10L 9/14****G10L 9/18**(21) Application number: **07295266**(22) Date of filing: **14 . 11 . 95**(30) Priority: **21 . 11 . 94 US 94 342494**(71) Applicant: **ROCKWELL INTERNATL CORP**(72) Inventor: **SU HUAN-YU  
LI TOM HONG**(54) **SYSTEM AND METHOD TO EVALUATE PITCH  
LAG AND VOICE ENCODER**

(57) Abstract:

**PROBLEM TO BE SOLVED:** To provide a compact and accurate pitch evaluating system into which a multiple resolution analysis to encode a voice is incorporated.

**SOLUTION:** The pitch evaluating device and method evaluate pitch lag of a voice by using a multiple resolution system. This system contains a step of sampling a voice, a step of alternately applying discrete Fourier transform and a step of squaring a result. Next, DTF is performed to convert a voice sample into a separate area to squared amplitude. Next, an initial pitch lag is obtained with low resolution. After the pitch lag is evaluated with low resolution, algorithm made precise is on the basis of minimizing an anticipating error in a time area. Next, a pitch lag made precise can be directly used in encoding a voice.

COPYRIGHT: (C)1996,JPO



BEST AVAILABLE COPY

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 08-211895

(43)Date of publication of application : 20.08.1996

(51)Int.Cl.

G10L 3/02

G10L 7/04

G10L 9/14

G10L 9/18

(21)Application number : 07-295266

(71)Applicant : ROCKWELL INTERNATL CORP

(22)Date of filing : 14.11.1995

(72)Inventor : SU HUAN-YU  
LI TOM HONG

(30)Priority

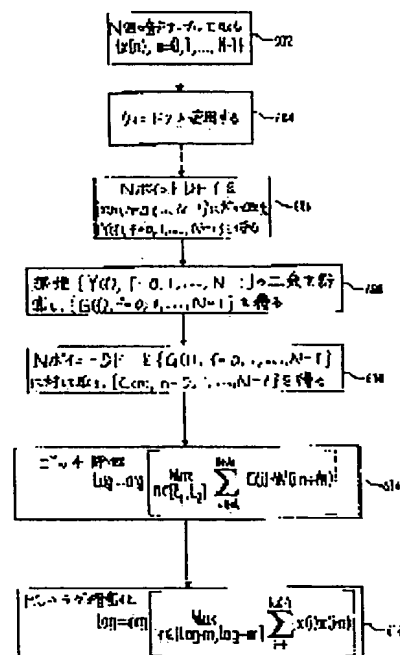
Priority number : 94 342494 Priority date : 21.11.1994 Priority country : US

## (54) SYSTEM AND METHOD TO EVALUATE PITCH LAG AND VOICE ENCODER

(57)Abstract:

PROBLEM TO BE SOLVED: To provide a compact and accurate pitch evaluating system into which a multiple resolution analysis to encode a voice is incorporated.

SOLUTION: The pitch evaluating device and method evaluate pitch lag of a voice by using a multiple resolution system. This system contains a step of sampling a voice, a step of alternately applying discrete Fourier transform and a step of squaring a result. Next, DTF is performed to convert a voice sample into a separate area to squared amplitude. Next, an initial pitch lag is obtained with low resolution. After the pitch lag is evaluated with low resolution, algorithm made precise is on the basis of minimizing an anticipating error in a time area. Next, a pitch lag made precise can be directly used in encoding a voice.



## LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision  
of rejection]

[Date of requesting appeal against examiner's  
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2000 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平8-211895

(43) 公開日 平成8年(1996)8月20日

(51) Int.Cl. <sup>8</sup>	識別記号	序内整理番号	F I	技術表示箇所
G 1 0 L 3/02	B			
7/04	G			
9/14	J			
9/18	E			

審査請求 未請求 請求項の数20 O L (全 13 頁)

(21) 出願番号 特願平7-295266

(22) 出願日 平成7年(1995)11月14日

(31) 優先権主張番号 08/342494

(32) 優先日 1994年11月21日

(33) 優先権主張国 米国 (US)

(71) 出願人 590002448

ロックウェル・インターナショナル・コー  
ポレーションROCKWELL INTERNATIONAL  
CORPORATIONアメリカ合衆国、90740-8250 カリフォル  
ニア州、シール・ビーチ、シールビー  
チ・プールバード、2201

(72) 発明者 フアン・ユー・スー

アメリカ合衆国、92714 カリフォルニア  
州、アーバイン、デュランゴ・アイル、  
398

(74) 代理人 弁理士 深見 久郎 (外3名)

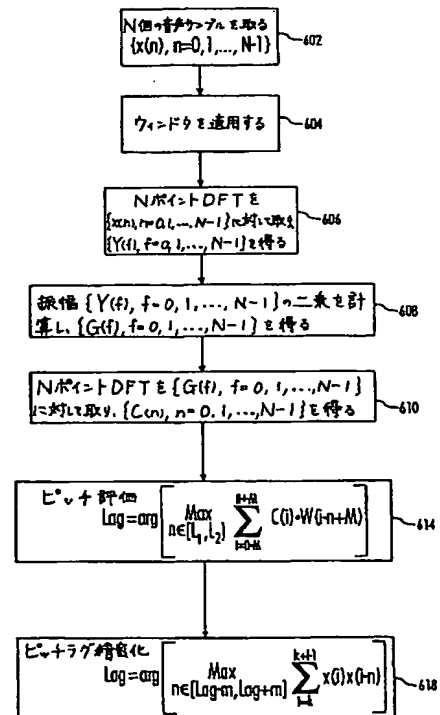
最終頁に続く

(54) 【発明の名称】 ピッチラグを評価するためのシステムおよび方法、ならびに音声符号化装置および方法

(57) 【要約】

【課題】 音声符号化のための多分解能分析を組入れる、簡潔かつ正確なピッチ評価システムを提供する。

【解決手段】 ピッチ評価装置および方法は多分解能方式を利用して音声のピッチラグを評価する。このシステムは音声をサンプルリングするステップと、交替に、離散フーリエ変換を適用するステップ、および結果を二乗するステップとを含む。次に、二乗された振幅に対して DFT が音声サンプルを別の域へ変換するために行なわれる。初期のピッチラグは次に、低分解能で求められ得る。低分解能ピッチラグ評価を得た後で、精密化されたアルゴリズムが適用されて高分解能ピッチラグを得る。精密化されたアルゴリズムは時間域で予測誤差を最小にすることに基づく。精密化されたピッチラグを次に、音声符号化において直接用いることができる。



## 【特許請求の範囲】

【請求項1】 音声量子化および圧縮のためにピッチラグを評価するためのシステムであって、音声は複数の音声サンプルによって規定され、現在の音声サンプルの評価は過去のサンプルの線形結合に従う時間域で決定され、前記システムは、

音声サンプルに第1の離散フーリエ変換(DFT)を適用するための手段を含み、前記第1のDFTは関連した振幅を有し、さらに、

前記第1のDFTの前記振幅を二乗するための手段と、前記二乗された振幅に対して第2のDFTを適用するための手段と、

時間域変換された音声サンプルに従って初期のピッチラグ値を決定するための手段と、

精密化されたピッチラグ値に従って前記音声サンプルを符号化するための手段とを含む、音声量子化および圧縮のためにピッチラグを評価するためのシステム。

【請求項2】 前記初期のピッチラグ値は関連した予測誤差を有し、前記システムは前記初期のピッチラグ値を精密化するための手段をさらに含み、前記関連した予測誤差は最小にされる、請求項1に記載のシステム。

【請求項3】 前記複数の音声サンプルを現在の符号化フレームへ分類するための手段と、

前記符号化フレームを複数のピッチサブフレームへ分割するための手段と、

前記ピッチサブフレームを複数の符号化サブフレームへ細分するための手段と、

前記現在の符号化フレームにおける各ピッチサブフレームの最後の符号化サブフレームに対して、それぞれラグ評価を表わす初期のピッチラグ評価  $lag_1$  および  $lag_2$  を評価するための手段と、

先行の符号化フレームにおける第2のピッチサブフレームのピッチラグ評価  $lag_0$  を精密化するための手段と、

$lag_1$ 、 $lag_2$ 、および  $lag_0$  を線形的に補間し、前記符号化サブフレームのピッチラグ値を評価するための手段と、

各符号化サブフレームの補間されたピッチラグをさらに精密化するための手段とをさらに含む、請求項1に記載のシステム。

【請求項4】 少数のサンプルで概略的に表わすためにダウンサンプリング値へ前記音声サンプルをダウンサンプリングするための手段をさらに含む、請求項1に記載のシステム。

【請求項5】 前記初期のピッチラグ値は式 ( $Lag_{scaled}$  = 音声サンプルの数 / ダウンサンプリング値) によって補正される、請求項4に記載のシステム。

【請求項6】 前記初期のピッチラグ値を精密化するための前記手段は自己相関を含む、請求項1に記載のシステム。

【請求項7】 前記音声サンプルを受けるための音声入力手段と、

前記精密化されたピッチラグ値を処理し、符号化された音声として入力音声を再現するためのコンピュータと、前記符号化された音声を出力するための音声出力手段とをさらに含む、請求項1に記載のシステム。

【請求項8】 入力音声を再現および符号化するための音声符号化装置であって、前記音声符号化装置は線形予測符号化(LPC)パラメータと、音声を発生するために音声再現を誘起するように参照される複数のベクトルを表わす新規コードブックとを用いるようにされており、前記音声符号化装置は、

前記入力音声を受けるための音声入力手段と、

前記入力音声を処理するためのコンピュータとを含み、前記コンピュータは、

前記入力音声内の現在の符号化フレームを切出すための手段と、

前記符号化フレームを複数のピッチサブフレームへ分割するための手段と、

N個の音声サンプルを有するピッチ分析ウィンドウを規定するための手段とを含み、前記ピッチ分析ウィンドウは前記ピッチサブフレームに対して延び、前記コンピュータは、

各ピッチサブフレームに対して初期のピッチラグ値を評価するための手段と、

各ピッチサブフレームを複数の符号化サブフレームへ分割するための手段とを含み、各ピッチサブフレームに対する前記初期のピッチラグ評価は、前記現在の符号化フレームにおける各ピッチサブフレームの最後の符号化サブフレームに対するラグ評価を表わし、前記コンピュータは、

前記評価されたピッチラグ値を前記ピッチサブフレームの間で線形的に補間し、各符号化サブフレームに対してピッチラグ評価を決定するための手段と、

各符号化サブフレームの前記線形的に補間されたラグ値を精密化するための手段とを含み、前記装置はさらに、前記精密化されたピッチラグ値に従って再現された音声を出力するための音声出力手段を含む、入力音声を再現および符号化するための音声符号化装置。

【請求項9】 前記コンピュータは、

少数のサンプルで表わすためにダウンサンプリング値Xへ前記N個の音声サンプルをダウンサンプリングするための手段と、

補正されたラグ値  $Lag_{scaled} = N / X$  であるように前記ピッチラグ値を補正するための手段とをさらに含む、請求項8に記載の装置。

【請求項10】 サンプリング速度Rで前記入力音声をサンプリングするサンプリング手段をさらに含み、前記N個の音声サンプルは式  $N = R * X$  によって決定される、請求項8に記載の装置。

【請求項 1 1】  $X=25\text{ms}$ 、 $R=8000\text{Hz}$ 、および  $N=320$  サンプルである、請求項 1 0 に記載の装置。

【請求項 1 2】 各符号化フレームはおおよそ  $40\text{ms}$  の長さを有する、請求項 8 に記載の装置。

【請求項 1 3】 音声量子化および圧縮のためにピッチラグを評価するための方法であって、前記音声は複数の音声サンプルによって規定され、現在の音声サンプルの評価は過去のサンプルの線形結合に従う時間域で決定され、前記方法は、

音声サンプルに第 1 の離散フーリエ変換 (DFT) を適用するステップを含み、前記第 1 の DFT は関連した振幅を有し、さらに、

前記第 1 の DFT の振幅を二乗するステップと、

前記第 1 の DFT の前記二乗された振幅に対して第 2 の DFT を適用するステップと、

時間域変換された音声サンプルに従って初期のピッチラグ値を決定するステップとを含み、前記初期のピッチラグ値は関連した予測誤差を有し、前記方法はさらに、自己相関を用いて前記初期のピッチラグ値を精密化するステップを含み、前記関連した予測誤差は最小にされ、さらに、

前記精密化されたピッチラグ値に従って前記音声サンプルを符号化するステップを含む、音声量子化および圧縮のためにピッチラグを評価するための方法。

【請求項 1 4】 前記複数の音声サンプルを現在の符号化フレームへ分類するステップと、

前記符号化フレームを複数のピッチサブフレームへ分割するステップと、

前記ピッチサブフレームを複数の符号化サブフレームへ細分するステップと、

前記現在の符号化フレームにおける各ピッチサブフレームの最後の符号化サブフレームに対して、それぞれラグ評価を表わす初期のピッチラグ評価  $lag_1$  および  $lag_2$  をそれぞれ評価するステップと、

先行の符号化フレームにおける第 2 のピッチサブフレームのピッチラグ評価  $lag_0$  を精密化するステップと、 $lag_1$ 、 $lag_2$ 、および  $lag_0$  を線形的に補間し、前記符号化サブフレームのピッチラグ値を評価するステップと、

各符号化サブフレームの補間されたピッチラグをさらに精密化するステップとをさらに含む、請求項 1 3 に記載の方法。

【請求項 1 5】 少数のサンプルで概略的に表わすためにダウンサンプリング値へ前記音声サンプルをダウンサンプリングするステップをさらに含む、請求項 1 3 に記載の方法。

【請求項 1 6】 式 ( $Lag_{scaled}$  = 音声サンプルの数 / ダウンサンプリング値) に従って前記初期のピッチラグ値を補正するステップをさらに含む、請求項 1 5 に記

載の方法。

【請求項 1 7】 前記音声サンプルを受けるステップと、前記精密化されたピッチラグ値を処理し、符号化された音声として前記入力音声を再現するステップと、前記符号化された音声を出力するステップとをさらに含む、請求項 1 3 に記載のシステム。

【請求項 1 8】 入力音声を再現および符合化するための音声符号化方法であって、音声符号化装置は線形予測符号化 (LPC) パラメータと、音声を発生するために音声再現を誘起するように参照される複数のベクトルを形成する擬似ランダム信号を表わす新規コードブックとを用いるようにされており、前記音声符号化方法は、前記入力音声を受取り、処理するステップと、前記入力音声を処理するステップとを含み、前記処理するステップは、

前記入力音声内で音声符号化フレームを決定するステップと、

前記符号化フレームを複数のピッチサブフレームへ細分するステップと、

$N$  個の音声サンプルを有するピッチ分析ウィンドウを規定するステップとを含み、前記ピッチ分析ウィンドウは前記ピッチサブフレームにわたって延び、前記処理するステップは、

各ピッチサブフレームに対して初期のピッチラグ値を概略的に評価するステップと、

各ピッチサブフレームに対する初期のピッチラグ評価が各ピッチサブフレームの最後の符号化サブフレームに対するラグ評価を表わすように、各ピッチサブフレームを複数の符号化サブフレームへ分割するステップと、

前記評価されたピッチラグ値を前記ピッチサブフレームの間で補間し、各符号化サブフレームに対してピッチラグ評価を決定するステップと、

線形的に補間されたラグ値を精密化するステップとを含み、前記方法はさらに、

精密化されたピッチラグ値に従って再現された音声を出力するステップを含む、入力音声を再現および符号化するための音声符号化方法。

【請求項 1 9】 前記処理するステップは、少数のサンプルで表わすためにダウンサンプリング値  $X$  へ前記  $N$  個の音声サンプルをダウンサンプリングするステップと、

補正されたラグ値  $Lag_{scaled} = N/X$  であるように前記ピッチラグ値を補正するステップとをさらに含む、請求項 1 8 に記載の装置。

【請求項 2 0】 前記  $N$  個の音声サンプルが式  $N = R * X$  に従って決定されるように、サンプリング速度  $R$  で前記入力音声をサンプリングするステップをさらに含む、請求項 1 8 に記載の方法。

【発明の詳細な説明】

【0001】

【発明の背景】信号のモデル化およびパラメータ評価はデータ圧縮、復元、および符号化においてますます重要な役割を果たす。基本的な話声をモデル化するために、音声信号は離散波形としてサンプリングされ、デジタル的に処理されなければならない。線形予測符号化(LPC)と称されるあるタイプの信号符号化技術において、何らかの特定の時間指標での信号値は前の値の線形関数としてモデル化される。こうして、後の信号はこれまでの値に従って線形的に予測される。結果として、効果的な信号表現は信号を表わすために、ある予測パラメータを評価し、かつ適用することによって決定できる。現在、符号励起線形予測(CELP)を含む音声符号化のためにLPC技術が用いられている。

【0002】ピッチ情報は符号化の目的に対して、確かな音の指標および表示であると認められている。ピッチは話者の音声の基本的な特徴またはパラメータを記述する。人間の音声は一般に容易には数学的に定量化できないので、音声のピッチデータを効果的に評価できる音声評価モデルが、よりの確かつ正確に符号化され、かつ復号された音声を提供する。しかしながら、あるCELP(たとえば、ベクトルと励起線形予測(VSELP)、マルチパルス、正規パルス、代数的CELPなど)およびMBEコーダ/デコーダ(「コーデック」)のような現在の音声符号化モデルにおいて、ピッチ評価アルゴリズムは正確さが高く、かつ複雑さが低いことを必要とするために、ピッチ評価がしばしば困難である。

【0003】いくつかのピッチラグ評価方式は上述されたコーデック(時間域方式、周波数域方式、およびケプ

$$A(z) = 1 - a_1 z^{-1} - a_2 z^{-2} - \dots - a_{np} z^{-np}$$

またはn番目のサンプルが次式によって予測できる

$$\hat{y}(n) = \sum_{k=1}^{np} a_k * y(n-k)$$

【0007】この式において、npはLPC予測次数(通例、約10)であり、y(n)はサンプリングされた音声データであり、nは時間指標を表わす。上記のLPCの式は、過去のサンプルの線形結合に従って現在のサンプルの評価を記述する。人間の耳の感度をモデルとするLPCフィルタに基づく聴感補正フィルタはここで次式によって規定される。

【0008】

【数2】

$$W(z) = \frac{A(z/\gamma_1)}{A(z/\gamma_2)} \quad \text{但し } 0 < \gamma_2 < \gamma_1 \leq 1$$

【0009】所望のピッチパラメータを抽出するために、次の重み付き符号化誤差エネルギーを最小にするピッチパラメータは各符号化サブフレームについて計算されなければならない。ここで1つの符号化フレームは、分析

ストラム域方式)と関連して用いられる。ピッチラグおよび音声再現の間に密接な関係があるために、ピッチ評価の正確さは通話品質に直接的な影響を及ぼす。CELPコーダでは、音声発生は予測(長期ピッチ予測および短期線形予測)に基づく。図1は典型的なCELPコーダによる音声再生のブロック図を示す。

【0004】音声データを圧縮するためには、不可欠な情報のみを抽出して冗長の伝送を回避することが望ましい。音声は短いブロックに分類でき、ここで代表的なパラメータはあらゆるブロックにおいて識別できる。図1に示されるように、CELP音声コーダは良質な音声を発生するために、LPCパラメータ110と、(ラグおよびその係数を含む)ピッチラグパラメータ112と、その利得パラメータ116を備える最適な新規コードベクトル114とを符号化されるべき入力音声から抽出しなければならない。コーダは適当な符号化方式を実現することによってLPCパラメータを量子化する。各パラメータの量子化の指標は、音声デコーダに記憶または伝送されるべき情報を含む。CELPコーデックでは、ピッチ予測パラメータ(ピッチラグおよびピッチ係数)の決定は時間域で行なわれるが、MBEコーデックでは、ピッチパラメータは周波数域で評価される。

【0005】LPC分析の後でCELPエンコーダは、(通例約10-40msで取られる)現在の音声符号化フレームのために適当なLPCフィルタ110を決定する。LPCフィルタは次式によって表わされる。

【0006】

【数1】

および符号化のためにいくつかの符号化サブフレームへ分割できる。

【0010】

【数3】

$$d = \|T - \beta P_{Lag} H - \alpha C_i H\|^2$$

【0011】この式において、Tは知覚的にフィルタされた入力信号を表わす目標信号であり、HはフィルタW(z)/A(z)のインパルス応答行列である。PLagはピッチラグ「Lag」と、所定のラグについて独自に規定された予測係数βとを有するピッチ予測寄与であり、Ciはコードブックにおける指標iおよびその対応する利得αに関連したコードブック寄与である。典型的には、人間の音声のピッチは2msから20msの間で異なる。したがって、音声は8KHzのサンプリング速



度でサンプリングされると、ピッチラグは概算で20サンプルから147サンプルに対応する。さらに、 $i$ は0および $N_c - 1$ の間の値を取り、ここで $N_c$ は新規コードブックのサイズである。

【0012】1タップピッチ予測子および1つの新規コードブックを想定する。しかしながら、典型的にピッチ予測子の一般的な形状は多タップ方式であり、新規コードブックの一般的な形状は多レベルベクトル量子化であるか、または、複数の新規コードブックを用いる。より詳細には、音声の符号化において、1タップピッチ予測子は現在の音声サンプルが1つの過去の音声サンプルによって予測できることを示し、一方多タップ予測子は現

$$\underset{Lag \in [L_1, L_2]}{Max} \frac{(TH^T P_{Lag}^T)^2}{\|P_{Lag} H\|^2}$$

【0015】この時間域方式は真のピッチラグを決定できるが、高いピッチ周波数を有する女性の音声には、式(1)によって求められるピッチラグは真のラグではなく、真のラグの倍数となり得る。この評価誤差を回避するために、評価誤差を訂正（たとえば、ラグの平滑化）する付加的なプロセスが必要であり、これはそれと引換えに不所望な複雑さを引き起こす。

【0016】しかしながら、このように過度に複雑であることは、時間域方式を用いる際の著しい欠点である。たとえば、整数のラグのみを用いてラグを決定するために、時間域方式は1秒当り300万回の動作（3MOP）を少なくとも必要とする。さらに、ピッチラグの平滑化および分数のピッチラグが用いられるならば、複雑さはほぼ4MOPであろう。実際には、容認可能な正確さでフルレンジのピッチラグ評価を実行するために、概算で1秒当り600万回のデジタル信号処理機械命令（6DSP MIP）が必要とされる。したがって、ピッチ評価は4から6のDSP MIPを必要とするとして一般に認められている。ピッチ評価の複雑さを減少できる方式は他にもあるが、そのような方式はしばしば質を犠牲にする。

【0017】正弦コーダの類で重要な要素であるMBEコーダでは、符号化パラメータは周波数域において抽出され、かつ量子化される。MBE音声モデルは図2から図4に示される。図2および図3に説明されるMBE音声エンコーダ／デコーダ（「ボゴダ」）では、基本周波数（またはピッチラグ）210、有声／無声決定212、およびスペクトルエンベロープ214は周波数域において入力音声から抽出される。パラメータは次に、記憶または伝送できるビットストリームへ量子化され、かつ符号化される。

【0018】MBEボコーダでは、良質な音声を達成するために基本周波数が高い正確さで評価されなければならない。基本周波数の評価は2段階で行なわれる。第1

在の音声サンプルが複数の過去の音声サンプルによって予測できることを意味する。

【0013】複雑さについて懸念があるために、最適な方式に準ずる方式が音声符号化方式において用いられてきた。たとえば、ピッチラグ評価は2.5msから18.5msをカバーするために、 $L_1$  および  $L_2$  サンプルの間の範囲で起こり得るラグ値を単に評価することによってなされてもよい。したがって、評価されたピッチラグ値は次式を最大にすることによって決定される。

【0014】

【数4】

式(1)

に、初期のピッチラグが21サンプルから114サンプルの範囲内で探索され、周波数域において入力音声216および合成された音声218の間で重み付き平均二乗誤差方程式310（図3）を最小にすることによって8000Hzのサンプリング速度で2.6msから14.25msをカバーする。元の音声および合成された音声の間の平均二乗誤差は次式によって与えられる。

【0019】

【数5】

$$E = \frac{1}{2\pi} \int_{-\pi}^{\pi} G(\omega) |s(\omega) - \hat{S}(\omega)| d\omega$$

【0020】この式において、 $S(\omega)$  は元の音声スペクトルであり、 $\hat{S}(\omega)$ （ $\hat{\cdot}$ は大文字Sの上にあると見なされる）は合成された音声スペクトルであり、 $G(\omega)$  は周波数依存重み付き関数である。図4に示されるように、ピッチ追跡アルゴリズム410は、隣接するフレームのピッチ情報を用いることによって、初期のピッチラグ評価412を更新するのに用いられる。

【0021】この方式を用いるのは、基本周波数が隣接するフレームの間で不意には変化するはずはないという仮定のためである。2つの過去の隣接するフレームおよび2つの未来の隣接するフレームのピッチ評価はピッチ追跡のために使用される。次に、（2つの過去のフレームおよび2つの未来のフレームを含む）平均二乗誤差は最小にされて現在のフレームの新しいピッチラグ値を求める。初期のピッチラグを追跡した後で、ピッチラグ多重検査方式414が多重ピッチラグを除去するために適用され、ピッチラグを平滑化する。

【0022】図4を参照すると、基本周波数評価の第2段階でピッチラグ精密化416が用いられてピッチ評価の正確さを高める。ピッチラグ候補値は初期のピッチラグ評価に基づいて形成される（すなわち、新しいピッチラグ候補値は、初期のピッチラグ評価からある分数を加

算し、または減算することによって形成される)。したがって、精密化されたピッチラグ評価418は、平均二乗誤差関数を最小にすることによってピッチラグ候補の中で決定できる。

【0023】しかしながら、周波数域ピッチ評価はある欠点を有する。第1に、複雑さが非常に高い。第2に、ピッチラグは2.5msから14.25msしかカバーしない20および114サンプルの範囲内で探索されて、256ポイントFFTに対処するように256サンプルにウィンドウサイズを制限しなければならない。しかしながら、非常に低いピッチ周波数の話者には、または14.25msを超えるピッチラグを有する音声には、256サンプルウィンドウ内で十分な数のサンプルを集めるのが不可能である。さらに、音声フレームにわたって評価されるのは、平均されたピッチラグだけである。

【0024】1967年にエイ・エム・ノル(A.M.Noll)によって提案されたケプストラム域ピッチラグ評価(図5)を用いて、変形された方法が他に提案された。ケプストラム域ピッチラグ評価では、510でおおよそ37msの音声サンプルがサンプリングされるので、可能な最大のピッチラグ(たとえば、18.5ms)の少なくとも2周期がカバーされる。次に、512ポイントFFTは音声フレームウィンドウに(ブロック512で)適用され、周波数スペクトルを獲得する。周波数スペクトルの対数514の振幅を取って、別の512ポイント逆FFT516がケプストラムを得るために適用される。重み付け関数518はケプストラムに適用され、ケプストラムのピークはピッチラグを決定するために520で検出される。次に、追跡アルゴリズム522が実行されて、いかなるピッチ倍数をも除去する。

【0025】しかしながら、ケプストラムピッチ検出方法にはいくつかの欠点が見受けられる。たとえば、計算上の要求が高い。8kHzのサンプリング速度において20サンプルおよび147サンプルの間でピッチの範囲をカバーするために、512ポイントFFTは二度行なわれなければならない。ケプストラムピッチ評価が分析フレームにわたる平均されたピッチラグの評価のみを提供するので、評価の正確さが不十分である。しかしながら、低ビット転送速度音声符号化については、ピッチラグ値が短い時間期間にわたって評価されることは重要である。結果として、ケプストラムピッチ評価は今日、高質な低ビット転送速度音声符号化についてはほとんど用いられない。したがって、上述された方式の各々に制限があるために、効果的なピッチラグ評価のための手段に

は、高質な低ビット転送速度音声符号化の必要を満たすことが所望される。

【0026】

【発明の概要】したがって、この発明の目的は、複雑さが最小であって正確さが高いことを必要とする、音声符号化のための多分解能分析を組入れるピッチ評価システムを提供することである。特定の実施例では、この発明はCELP技術ならびにさまざまな他の音声符号化および認識システムを用いる音声符号化の装置および方法を対象とする。したがって、必要な高い正確さを維持しながら、よりよい結果がより少ない計算手段でもたらされる。

【0027】これらの目的および他の目的は、この発明の実施例に従って、音声の的確な再現および再生を速くかつ効果的に可能にするピッチラグ評価方式によって達成される。ピッチラグは所定の音声フレームについて抽出され、次に、各サブフレームについて精密化される。最小の数の音声サンプルが音声を直接サンプリングすることによって獲得された後で、離散フーリエ変換(DFT)が適用され、結果として生じる振幅が二乗される。第2のDFTが次に行なわれる。したがって、フレーム内の音声サンプルに対する的確な初期のピッチラグは、8kHzのサンプリング速度で20サンプルの可能な最小値と147サンプルの最大ラグ値との間で決定できる。初期のピッチラグ評価を獲得した後で、時間域精密化がさらに評価の正確さを向上するために各サブフレームについて行なわれなければならない。

【0028】

【好ましい実施例の詳細な説明】この発明の好ましい実施例に従ったピッチラグ評価方式は一般に図6、7、8および9において示される。まず、N個の音声サンプル $\{x(n), n=0, 1, \dots, N-1\}$ が集められる。

(図6のステップ602) Nはたとえば、8000Hzのサンプリング速度で典型的な40msの音声ウィンドウに対処するために320個の音声サンプルに等しくてもよい。Nの値はおおまかに評価された音声周期によって決定され、ここで少なくとも2周期が音声スペクトルを発生するために一般に必要とされる。このように、Nが可能な最大のピッチラグの2倍よりも大きくなくてはならず、ここでは $\{x(n), n=0, 1, \dots, N-1\}$ である。さらに、少なくとも2ピッチ周期をカバーするハミングウィンドウ604または他のウィンドウが好ましくは実現される。

【0029】

【数6】

NポイントDFTがステップ606で $\{x(n), n=0, 1, \dots, N-1\}$ にわたって適用され、振幅 $\{Y(f), f=0, 1, \dots, N-1\}$ を得る。ここで

$$Y(f) = \sum_{n=0}^{N-1} x(n)e^{-j2\pi nf/N} \quad \text{但し } f=0, 1, \dots, N-1 \quad \text{式(2)}$$

次に、 $Y(f)$ はステップ608で次式に従って二乗される。

$$G(f) = |Y(f)|^2 \quad \text{但し } f=0, 1, \dots, N-1 \quad \text{式(3)}$$

第2のNポイントDFTがステップ610で $G(f)$ に適用されて以下を獲得する。

$$C(n) = \sum_{f=0}^{N-1} G(f)e^{-j2\pi nf/N} \quad \text{但し } n=0, 1, \dots, N-1 \quad \text{式(4)}$$

【0030】この発明の実施例に従って、関数 $G(f)$ ではなく $G(f)$ の対数が式(4)において用いられる従来のケプストラム変換と $C(n)$ とが異なることが認識される。この違いの原因は一般的には複雑さである。除去されなければ実質的により多くの計算資源を必要とする対数関数を除去することによって複雑さを減少することが望ましい。さらに、ケプストラムまたは $C(n)$ 関数を用いるピッチラグ評価方式を比較すると、音声の無声または遷移区間のみに対して異なった結果が獲得されていた。たとえば、無声または遷移音声に対してピッチ

チの定義が不明確である。遷移音声にはピッチがないといわれてきたが、誤差を最小にするために何らかの予測を常に示すことができるともいわれる。

【0031】したがって、一旦 $C(n)$ が決定されると(ステップ610)、所定の音声フレームに対するピッチラグは次式を解くことによってステップ614で求められ得る。

【0032】

【数7】

$$Lag = \arg \left[ \max_{n \in [L_1, L_2]} \sum_{i=n-M}^{n+M} C(i) \cdot W(i-n+M) \right] \quad \text{式(5)}$$

【0033】この式において、 $\arg[\cdot]$ は内部最適化関数を満たす変数 $n$ を決定し、 $L_1$ および $L_2$ はそれぞれ可能な最小のピッチラグおよび可能な最大のピッチラグとして規定される。音声符号化の便宜上、 $L_2$ および $L_1$ の間の差は2進数表現のために2の累乗であることが望ましい。好ましい実施例では、 $L_1$ および $L_2$ はそれぞれ20および147の値を取って典型的な人間の音声のピッチラグ範囲の2.5msから18.375msをカバーし、ここで $L_1$ および $L_2$ の間の間隔は2の累乗である。 $W(i)$ は重み付き関数であり、 $2M+1$ はウィンドウサイズを表す。好ましくは、 $\{W(i) = 1, i=0, 1, \dots, 2M\}$ であり、 $M=1$ である。

【0034】結果として生じるピッチラグは平均された値であるが、それは信頼できて的確であるということがわかった。平均化から生じる効果は絶対的に大きな分析ウィンドウサイズによるものであり、147サンプルのラグに対して、ウィンドウサイズはラグ値の少なくとも2倍であるべきである。しかしながら、不所望なことに、典型的に小さいピッチラグを示す女性の話者のようなある話者からの信号は、このような大きなウィンドウ

では4から10ピッチ周期を含み得る。ピッチラグに変化があれば、提案されたピッチラグ評価は平均されたピッチラグしか生成しない。結果として、音声符号化においてこのような平均されたピッチラグを用いることで音声評価および再生に大きな劣化が生じ得る。

【0035】音声におけるピッチ情報の相対的に速い変化のために、CELPモデルに基づくほとんどの音声符号化システムはサブフレームごとに一度ピッチラグを評価し、かつ伝送する。こうして、典型的には2msから10msの長さ(16から80サンプル)であるいくつかの音声サブフレームへ1つの音声フレームが分割されるCELP型音声符号化において、ピッチ情報は各サブフレームで更新される。したがって、正確なピッチラグ値はサブフレームのためにのみ必要とされる。しかしながら、上記の方式に従って評価されたピッチラグは、平均化から生じる影響のために正確な音声符号化には十分な正確さを有さない。

【0036】こうして、この発明の特定の実施例において、評価の正確さを向上させるために、初期のピッチラグ評価に基づいた精密化された探索が時間域において行

なわれる（ステップ618）。簡単な自己相関方法がほぼ平均されたLag値で特定の符号化周期またはサブフレームに対して行なわれる。

$$Lag = \arg \left[ \underset{n \in [Lag-m, Lag+m]}{Max} \sum_{i=k}^{k+l-1} x(i)x(i-n) \right] \quad \text{式(6)}$$

【0038】この式において、 $\arg[\cdot]$ は内部最適化関数を満たす変数 $n$ を決定し、 $k$ はサブフレームの第1のサンプルを示し、 $l$ は精密化ウィンドウサイズを表わし、 $m$ は探索範囲である。的確なピッチラグ値を決定するために、精密化ウィンドウサイズは少なくとも1ピッチ周期であるべきである。しかしながら、ウィンドウは平均化の影響を避けるためにあまりに大きすぎてはならない。たとえば、好ましくは $l = Lag + 10$ 、および $m = 5$ である。こうして、式(6)の時間域精密化に従って、より正確なピッチラグが評価されてサブフレームの符号化に適用できる。

【0039】動作時において、高速フーリエ変換（FFT）が一般的なDFTよりも計算上効果的である場合もあるが、FFTを用いる際の欠点はウィンドウサイズが2の累乗でなければならないことである。たとえば、147サンプルの最大のピッチラグは2の累乗ではないことが示されてきた。最大のピッチラグを含むためには、512サンプルのウィンドウサイズが必要である。しかしながら、このことで、上述された平均化から生じる影

$$y(i) = x([i \cdot \lambda]) + \{x([i \cdot \lambda] + 1) - x([i \cdot \lambda])\}(i \cdot \lambda - [i \cdot \lambda])$$

【0037】

【数8】

響のために女性の音声に対するピッチラグ評価の質が悪くなり、多量の計算が必要となる。256サンプルのウィンドウサイズが用いられるならば、平均化から生じる影響は減少され、複雑さが一層少なくなる。しかしながら、このようなウィンドウを用いると音声中の128サンプルよりも大きなピッチラグには対処できない。

【0040】これらの問題のいくつかを克服するために、この発明の代替の好ましい実施例は256ポイントFFTを利用して複雑さを減少し、変更された信号を用いてピッチラグを評価する。信号を変更するのはダウンサンプリングプロセスである。図7および図8を参照すると、 $N$ 個の音声サンプル $\{x(n), n=0, 1, \dots, N-1\}$ が集められ（ステップ702）、 $N$ は最大のピッチラグの2倍よりも大きい。次に、 $N$ 個の音声サンプルが次式に従って、線形補間を用いて256個の新しい分析サンプルへダウンサンプリングされる（ステップ704）。

【0041】

【数9】

但し  $i = 0, 1, \dots, 255$

【0042】この式において、 $\lambda = N/256$ であり、角括弧内の値すなわち $[i \cdot \lambda]$ は $i \cdot \lambda$ 以下の最大の整数値を示す。次に、ステップ705でハミングウィンドウまたは他のウィンドウが補間されたデータに適用される。

【0043】ステップ706では、ピッチラグ評価は256ポイントFFTを用いて $y(i)$ にわたって行なわれ、振幅 $Y(f)$ を発生する。次に、ステップ708からステップ710は図6に関して説明されたのと同様に実行される。しかしながら、 $G(f)$ はさらにフィルタされ（ステップ709）、ピッチ検出のためには有用ではない、 $G(f)$ の高周波数成分を減少する。一旦 $y(i)$ のラグすなわち $Lag_y$ が式(5)に従って求められれば（ステップ714）、これはステップ716で再スケールされてピッチラグ評価を決定する。

【0044】

【数10】

$$Lag = Lag_y \cdot \lambda$$

【0045】要約すると、符号化フレームのための初期のピッチ評価を求める上記の手順は以下のとおりである。

【0046】(1) 標準40msの符号化フレームを

ピッチサブフレーム802および804へ細分する。各ピッチサブフレームはおおよそ20msの長さである。

【0047】(2) ピッチ分析ウィンドウ806が最後のサブフレームの中心に位置決めされるように $N=320$ 個の音声サンプルを取り、提案されたアルゴリズムを用いてそのサブフレームに対するラグを求める。

【0048】(3) ピッチサブフレームに対する初期のピッチラグ値を決定する。

次に、時間域精密化が元の音声サンプル $x(n)$ にわたってステップ718で行なわれる。こうして、この発明の実施例において、複雑さを減少してなお、高い正確さを維持しながらピッチラグ値が的確に評価できる。この発明のFFT実施例を用いると、120よりも大きいピッチラグ値に対処するのは困難ではない。

【0049】より詳細には、時間域精密化は元の音声サンプルにわたって行なわれる。たとえば、40msの符号化フレームは図9に示されるようにまず、8個の5msのサブフレーム808へ分割される。初期のピッチラグ評価 $lag_1$ および $lag_2$ は、現在の符号化フレームにおける各ピッチサブフレームの最後の符号化サブフレームに対するラグ評価である。 $lag_0$ は先行の符号化フレームにおける第2のピッチサブフレームの精密化

されたラグ評価である。 $lag_1$ 、 $lag_2$ 、および $lag_0$ の間の関係は図9に示される。

【0050】初期のピッチラグ $lag_1$ および $lag_2$ は次式に従って最初に精密化されて、その正確さを向上

$$lag_i = \arg \left[ \max_{n \in [lag_i - M, lag_i + M]} \sum_{k=N_i}^{N_i+L-1} x(k) \cdot x(k-n) \right]$$

但し  $i = 1, 2$

【0052】ここで $N_i$ は、ピッチ $lag_i$ に対するピッチサブフレームにおける開始サンプルの指標である。好ましくは、 $M$ は10と選択され、 $L$ は $lag_i + 10$ であり、 $i$ はピッチサブフレームの指標を示す。

【0053】一旦初期のピッチラグの精密化が完了すると、符号化サブフレームのピッチラグが決定できる。符号化サブフレームのピッチラグは $lag_1$ 、 $lag_2$ 、および $lag_0$ を線形的に補間することによって評価される。符号化サブフレームのピッチラグ評価の正確さ

$$lag_I(i) = \begin{cases} lag_0 + (lag_1 - lag_0) \cdot \frac{i+1}{4} & \text{但し } i = 0, 1, 2, 3 \\ lag_1 + (lag_2 - lag_1) \cdot \frac{i-3}{4} & \text{但し } i = 4, 5, 6, 7 \end{cases}$$

【0055】線形補間によって与えられるピッチラグ評価の正確さが十分ではないので、さらなる改良が必要となるであろう。与えられたピッチラグ評価 $\{lag_I(i), i = 0, 1, \dots, 7\}$ に対して、各 $lag_I$

$$lag(i) = \arg \left[ \max_{n \in [lag_I(i) - M, lag_I(i) + M]} \sum_{k=N_i}^{N_i+L-1} x(k) \cdot x(k-n) \right]$$

但し  $i = 0, 1, \dots, 7$

【0057】ここで $N_i$ はピッチ $lag(i)$ に対する符号化サブフレームにおける開始サンプルの指標である。例では、 $M$ は3と選択され、 $L$ は40に等しい。

【0058】さらに、ピッチラグの線形補間は音声の無声区間において重要である。何らかの分析方法によって求められたピッチラグは無声音声に任意に配分される傾向を有する。しかしながら、相対的に大きいピッチサブフレームサイズのために、各サブフレームに対するラグが(上の手順(2)で求められる)始めに決定されたサブフレームラグにあまりにも近い場合、元々は音声にはなかった不所望な人工の周期性が加えられる。さらに線形補間は、質の悪い無声音声に関連した問題を簡単に解決する。さらに、サブフレームのラグは任意である傾向を有するので、各サブフレームに対するラグは一旦補間されると、これも非常に任意に配分され、このことが音声の質を保証する。

【図面の簡単な説明】

【図1】CELP音声モデルのブロック図である。

させる(図8のステップ718)。

【0051】

【数11】

は、次の手順に従って各符号化サブフレームの補間されたピッチラグを精密化することによって向上する。精密化された初期のピッチ評価 $lag_1$ 、 $lag_2$ 、および $lag_0$ に基づく符号化サブフレームの補間されたピッチラグを $\{lag_I(i), i = 0, 1, \dots, 7\}$ が表わす場合、 $lag_I(i)$ は次式によって決定される。

【0054】

【数12】

( $i$ )は次式によってさらに精密化される(ステップ722)。

【0056】

【数13】

【図2】MBE音声モデルのブロック図である。

【図3】MBEエンコーダのブロック図である。

【図4】MBEボコーダにおけるピッチラグ評価のブロック図である。

【図5】ケプストラムに基づくピッチラグ検出方式のブロック図である。

【図6】この発明の実施例に従うピッチラグ評価の動作上のフロー図である。

【図7】この発明の別の実施例に従うピッチラグ評価のフロー図である。

【図8】この発明の別の実施例に従うピッチラグ評価のフロー図である。

【図9】図6の実施例に従う音声符号化の図である。

【符号の説明】

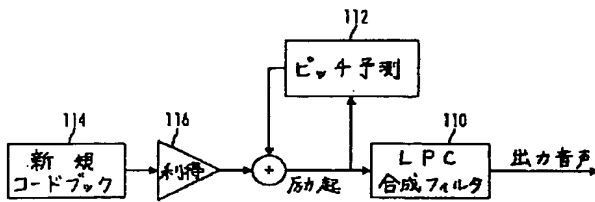
802 ピッチサブフレーム

804 ピッチサブフレーム

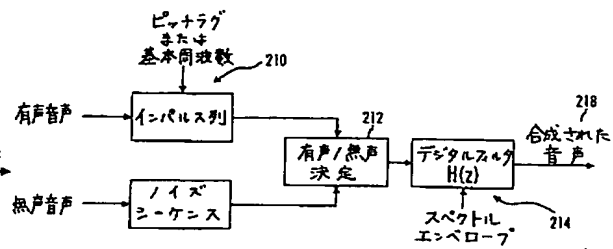
806 ピッチ分析ウィンドウ

808 サブフレーム

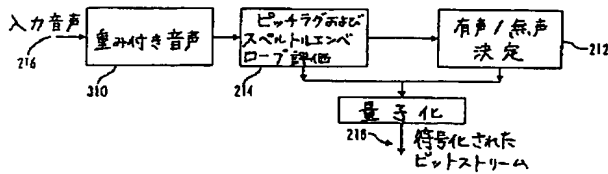
【図 1】



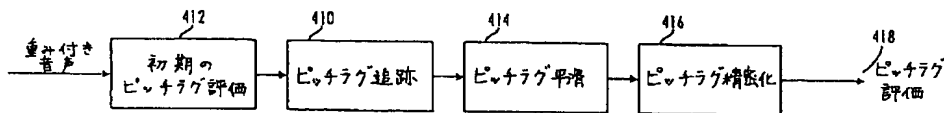
【図 2】



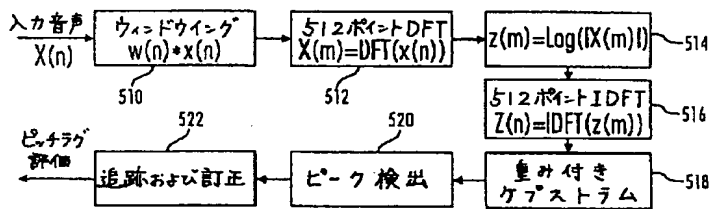
【図 3】



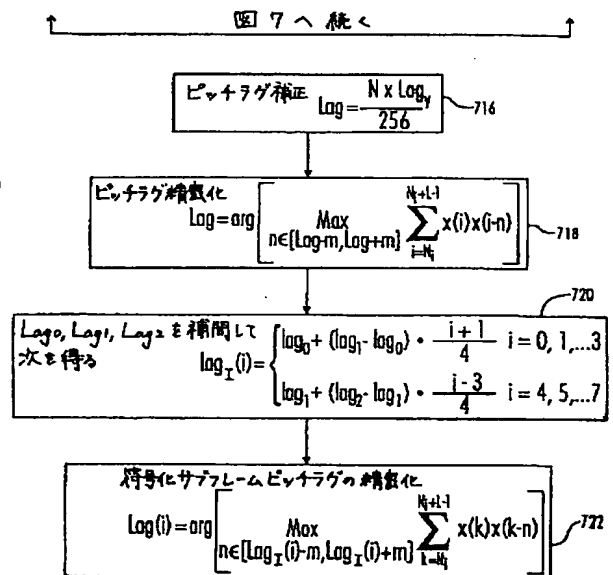
【図 4】



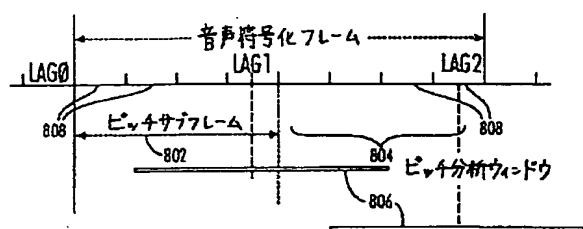
【図 5】



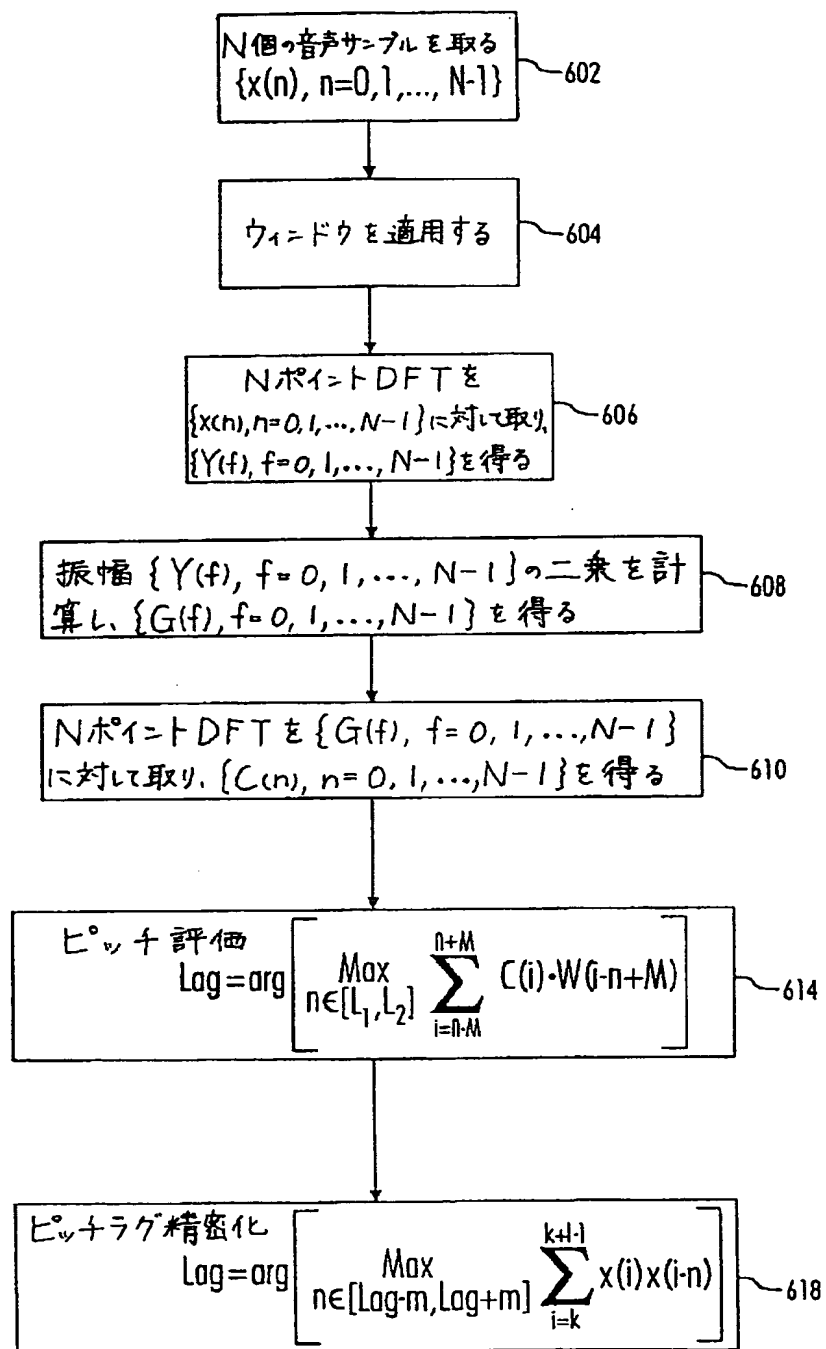
【図 8】



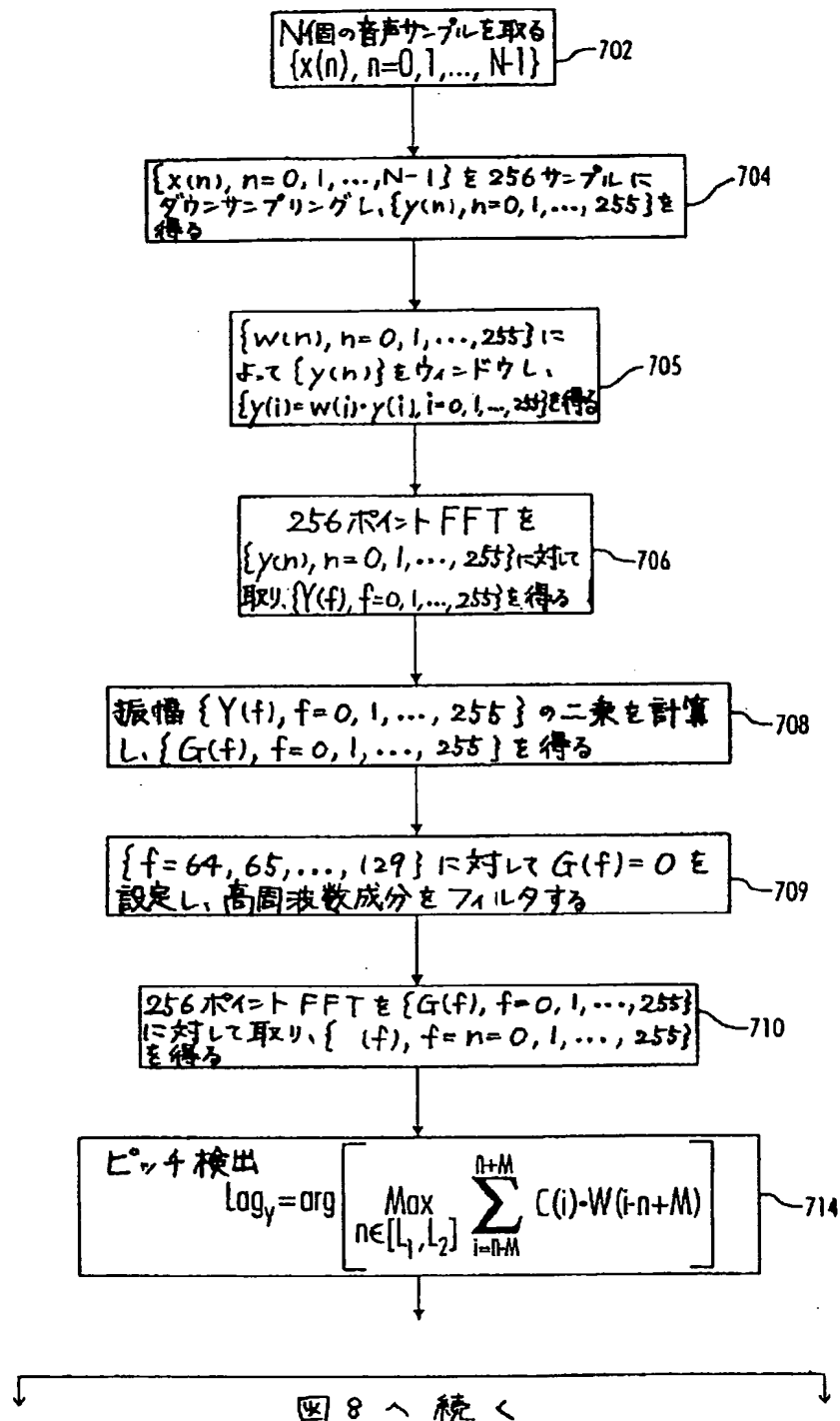
【図 9】



【図6】



【図7】





フロントページの続き

(72)発明者 トム・ホン・リー  
アメリカ合衆国、07748 ニュージャージー  
州、ミドルタウン、ノウルウッド・ドラ  
イブ、501